

Hans Reijnierse · Peter Borm  
Mark Voorneveld

## On ‘informationally robust equilibria’ for bimatrix games

Received: 5 October 2004 / Accepted: 7 December 2005 / Published online: 13 January 2006  
© Springer-Verlag 2006

**Abstract** Informationally robust equilibria (IRE) are introduced in Robson (*Games Econ Behav* 7:233–245, 1994) as a refinement of Nash equilibria for strategic games. Such equilibria are limits of a sequence of (subgame perfect) Nash equilibria in perturbed games where with small probability information about the strategic behavior is revealed to other players (information leakage). Focusing on bimatrix games, we consider a type of informationally robust equilibria and derive a number of properties: they form a non-empty and closed subset of the Nash equilibria. Moreover, IRE is a strict concept in the sense that the IRE are independent of the exact sequence of probabilities with which information is leaked. The set of IRE, like the set of Nash equilibria, is the finite union of polytopes. In potential games, there is an IRE in pure strategies. In zero-sum games, the set of IRE has a product structure and its elements can be computed efficiently by using linear programming. We also discuss extensions to games with infinite strategy spaces and more than two players.

**Keywords** Bimatrix game · Equilibrium selection · Leakage of information

**JEL Classification Numbers** C72

---

The authors would like to thank Marieke Quant for her helpful comments.

---

H. Reijnierse (✉) · P. Borm · M. Voorneveld  
Center and Department of Econometrics and Operations Research, Tilburg University,  
P.O. Box 90153, 5000 LE Tilburg, The Netherlands  
E-mail: J.H.Reijnierse@uvt.nl

M. Voorneveld  
Department of Economics, Stockholm School of Economics,  
Box 6501, 113 83 Stockholm, Sweden

## 1 Introduction

A branch of game-theoretic literature deals with refinements of the Nash equilibrium concept. Starting with the perfect equilibria of Selten (1975), and along the way developing notions like properness (Myerson 1978) and strict perfectness (Okada 1984), it eventually culminated in the work of (Kohlberg and Mertens 1986). See Van Damme (1991) for an overview. The original idea underlying these concepts is that players undergo a thought experiment in which all players make mistakes with small, but positive probabilities. The current paper uses a similar, but somewhat different thought experiment of the type suggested by Robson (1994), where with small, positive probability one of the players' action choices is revealed. Such "information leakage" is of relevance in numerous practical situations, witnessing the literature on industrial espionage, creating first- or second-mover (dis)advantages (see, for instance, Bagwell 1995), enforcing cooperation (Matsui 1989), but also – more casually – the importance of being able to hide the strength of your hand in a poker game.

Throughout the paper, we consider mixed extensions of finite, two-person games (bimatrix games). Games are perturbed by allocating small probabilities to two disjoint events. With large probability, the original game is played, but there is a small probability that the action choice of one of the players is revealed to the other. Informally, there is a small probability that one of the players acts first. If, say, player 1 acts first, player 2 observes the decision of player 1. If player 1 plays a mixed strategy, player 2 is informed about the outcome of the chance mechanism. Thereafter, he can base his decision on this information. Player 1 cannot distinguish between this case and the regular one, i.e., he does not know if he is revealing his action or not. Similarly, player 2 may act first (not knowing this himself) and player 1 can respond. The events *player 1 acts first* and *player 2 acts first* do not necessarily have the same probability.

An alternative way to model leakage of information is given in Reny and Robson (2004). Here, a player's *mixed* strategy may be revealed and the corresponding perturbed games are Bayesian. In particular, this allows for a reinterpretation of a mixed strategy Nash equilibrium by means of Bayesian equilibria. Solan and Yariv (2004) provide a two-player model in which player 1 can purchase (obtain by espionage) a noisy signal of the chosen strategy of player 2.

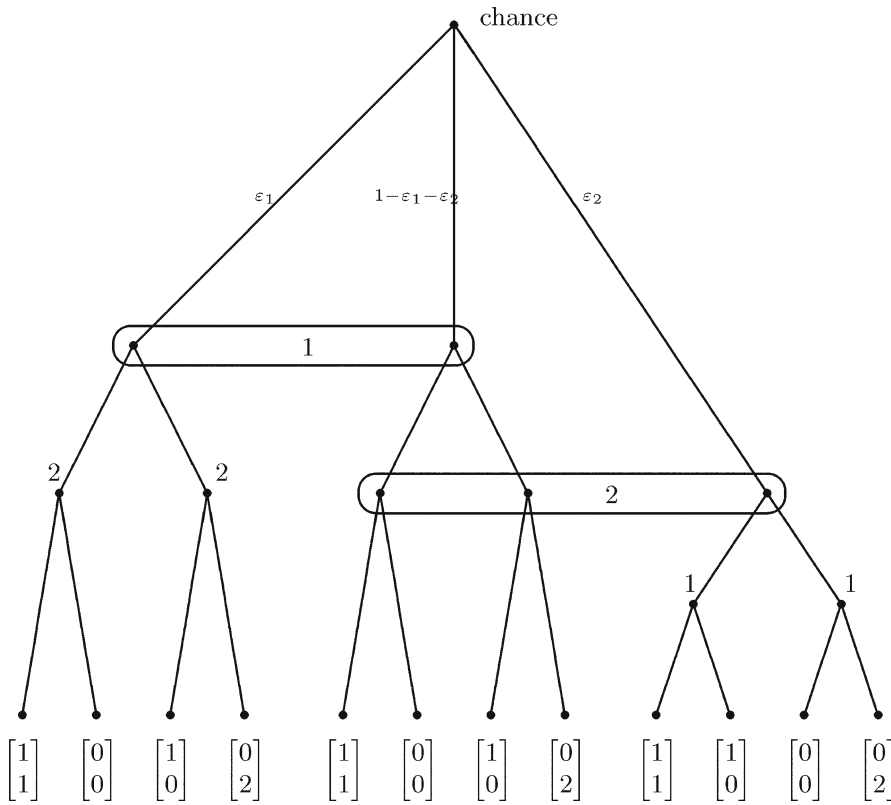
Our approach, however, follows the lines set out by Robson (1994). Focussing on bimatrix games, our underlying thought experiment allows for the possibility of a player's action to be revealed. Further, we put restrictions on the perturbations in order to have perturbed games of the same dimension as the original one.

To highlight these differences and to get acquainted with the model, let us discuss an example. Consider the bimatrix game

$$\begin{bmatrix} (1, 1) & (0, 0) \\ (1, 0) & (0, 2) \end{bmatrix}.$$

The row player has no direct influence on his payoff by his own action. He can however, have the following line of thought:

If there would be a slight chance that my opponent can act upon my action, then I'd better play the top row; my opponent's best reply to this action is playing the left column. This leads to a benefit of 1.



**Fig. 1** A perturbation concerning leakage of information in extensive form

Let  $\varepsilon_i$  denote the probability that player  $i$ 's action is revealed to the other player ( $i \in \{1, 2\}$ ). The extensive form of this perturbation is depicted in Figure 1. The analysis of Robson takes place in this extensive form game, requiring subgame perfection. Notice that there is an exponential growth in the size of strategy spaces in the perturbed game; if the players in the bimatrix game have  $m$ , respectively  $n$ , pure strategies, they have  $m^{n+1}$  and  $n^{m+1}$  pure strategies in the perturbed game. Our paper contains an analysis of perturbed games in normal form and avoids the exponential growth by putting a common rationality restriction on the behavior of the players. These restrictions only have a bite in the one-person (like) perfect information subgames at the end of the game tree. In a subgame perfect equilibrium, the player involved chooses an action that maximizes his utility. To reduce the strategy spaces of the perturbed games, we delete all other strategies *beforehand*. Generically, all but one strategy in such subgames are omitted.

As a tie-breaking rule in non-generic cases, we assume that the player involved chooses a utility maximizing action that maximizes the payoff to the other player ("optimistic tie-breaking").<sup>1</sup>

<sup>1</sup> An alternative tie-breaking rule would be to choose a utility maximizing action that is worst for the opponent (pessimistic tie-breaking). This would not affect the results of this paper. We thank a referee for pointing out that the optimistic tie is superior for reasons of generalization. For a discussion we refer to Section 7.

Having defined the perturbed games in this way, a strategy profile is called informationally robust if it is the limit of a sequence of Nash equilibria of perturbed games, with probabilities of information leakage converging to zero.

Our main results are as follows. Informationally robust equilibria are defined in Section 2. Furthermore, it provides an alternative way to describe informationally robustness (Lemma 1) and shows that the set of informationally robust equilibria is a non-empty, closed subset of the game's Nash equilibria (Theorem 1). Section 3 provides a second characterization of IRE by showing that the exact sequence of probabilities with which information leakage occurs is irrelevant (Theorem 2). Theorem 3 in Section 4 characterizes the structure of the set of IRE: like the set of Nash equilibria of bimatrix games it is the finite union of polytopes. Next, we consider two special classes of games. In Section 5 it is shown that in potential games (cf. Monderer and Shapley 1996) there is always a pure-strategy equilibrium that is informationally robust (Theorem 4). In Section 6 it is shown that – again in correspondence with the set of Nash equilibria, which is the Cartesian product of a zero-sum game's maximin/minimax strategies – the set of IRE in two-person zero-sum games has a product structure, whose elements can be computed efficiently using linear programming (Theorem 5). Section 7 concludes with a discussion of the possibilities of generalizing our analysis to games with infinite strategy spaces and games with more than two players.

## 2 IRE

Let us fix the notations that are used throughout the paper. A bimatrix game is the mixed extension of a finite two-person noncooperative game. It is characterized by a pair  $(A, B)$  of real-valued matrices of equal, finite, size. The players are called 1 and 2. Player 1 chooses a row and player 2 chooses a column. We use  $m$  for the number of rows and  $n$  for the number of columns. The index sets of the rows and columns are denoted by  $M$  and  $N$ , respectively:

$$M = \{1, \dots, m\} \text{ and } N = \{1, \dots, n\}.$$

Typical characters to index rows are  $i$  and  $k$ , typical characters to index columns are  $j$  and  $\ell$ . The spaces of mixed strategies are called  $\Delta_m$  and  $\Delta_n$ , respectively. Furthermore,  $\Delta = \Delta_m \times \Delta_n$ ; the space of strategy profiles. The unit vectors of  $\Delta_m$  and  $\Delta_n$  (i.e., the pure strategies) are denoted by  $e_i$  ( $i \in M$ ) and  $f_j$  ( $j \in N$ ). A typical element of  $\Delta_m$  will be denoted by  $p$ , a typical element of  $\Delta_n$  by  $q$ . Players have a pure best reply correspondence:

$$PB_1(A, q) = \operatorname{argmax}_{i \in M} e_i A q \text{ and } PB_2(B, p) = \operatorname{argmax}_{j \in N} p B f_j.$$

These correspondences are upper semi-continuous in both coordinates, e.g., if  $(A_t, q_t)$  tends to  $(A, q)$ , then  $PB_1(A_t, q_t) \subseteq PB_1(A, q)$  for sufficiently large  $t$ . The *carrier*  $C(x)$  of a vector  $x$  is the set of its non-zero coordinates, i.e.,

$$C(x) = \{i \mid x_i \neq 0\}.$$

A Nash equilibrium  $(p, q)$  is a profile of mixed strategies such that  $C(p) \subseteq PB_1(A, q)$  and  $C(q) \subseteq PB_2(B, p)$ . The set of all Nash equilibria of the game

$(A, B)$  is denoted by  $E(A, B)$ . Two extra parameters are needed to give the perturbations of  $A$  and  $B$ . The probability that the action of player 1 is revealed to player 2 is called  $\varepsilon_1 > 0$ . The probability that the action of player 2 is revealed to player 1 is called  $\varepsilon_2 > 0$ . By assumption,  $\varepsilon_1 + \varepsilon_2 < 1$ . We define  $A_{ij}(\varepsilon_1, \varepsilon_2)$ , the payoff player 1 receives in a perturbed game when player 1 chooses strategy  $e_i$  and player 2 chooses  $f_j$ , as follows. With large probability  $(1 - \varepsilon_1 - \varepsilon_2)$  he receives the original payoff  $A_{ij}$ . With probability  $\varepsilon_1$ , first player's action  $e_i$  is revealed to player 2, who can respond optimally to it, i.e., choose an element of  $PB_2(B, e_i)$ . In case of multiple best replies player 2 selects one of the strategies  $f_\ell \in PB_2(B, e_i)$  that maximizes his opponent's utility  $A_{i\ell}$ . Conversely, with probability  $\varepsilon_2$ , second player's action  $f_j$  is revealed to player 1, who reacts optimally against it, resulting in  $\max_{k \in M} A_{kj}$ . The perturbed game for player 2 is defined analogously. This leads to the following definition.

**Definition 1** Let  $(A, B)$  be an  $m \times n$ -bimatrix game and let  $(\varepsilon_1, \varepsilon_2)$  be a pair of positive real numbers satisfying  $\varepsilon_1 + \varepsilon_2 < 1$ . The perturbed game  $(A(\varepsilon_1, \varepsilon_2), B(\varepsilon_1, \varepsilon_2))$  is the bimatrix game given by

$$\begin{aligned} A_{ij}(\varepsilon_1, \varepsilon_2) &= (1 - \varepsilon_1 - \varepsilon_2)A_{ij} + \varepsilon_1 \max_{\ell \in PB_2(B, e_i)} \{A_{i\ell}\} + \varepsilon_2 \max_{k \in M} A_{kj}, \\ B_{ij}(\varepsilon_1, \varepsilon_2) &= (1 - \varepsilon_1 - \varepsilon_2)B_{ij} + \varepsilon_1 \max_{\ell \in N} B_{i\ell} + \varepsilon_2 \max_{k \in PB_1(A, f_j)} \{B_{kj}\}. \end{aligned}$$

Now we have made all preparations to define informationally robust equilibria.

**Definition 2** Let  $(A, B)$  be an  $m \times n$ -bimatrix game. A profile  $(p, q)$  is an informationally robust equilibrium or IRE if there exist sequences  $(\varepsilon_1^t)_{t \in \mathbb{N}}$  and  $(\varepsilon_2^t)_{t \in \mathbb{N}}$  of positive real numbers converging to zero, and a sequence  $(p^t, q^t)_{t \in \mathbb{N}}$  in  $\Delta$  converging to  $(p, q)$  such that for all  $t \in \mathbb{N}$ ,

$$(p^t, q^t) \in E(A(\varepsilon_1^t, \varepsilon_2^t), B(\varepsilon_1^t, \varepsilon_2^t)).$$

The set of informationally robust equilibria of  $(A, B)$  is denoted by  $IRE(A, B)$ .

There is an alternative convenient characterization of IRE by means of best reply equivalent perturbed games. Two bimatrix games  $(A, B)$  and  $(C, D)$  of equal size are called *best reply equivalent* if their pure best reply functions coincide:

$$PB_1(A, \cdot) = PB_1(C, \cdot) \text{ and } PB_2(B, \cdot) = PB_2(D, \cdot).$$

We will denote this type of equivalence by  $(A, B) \equiv_b (C, D)$ .

Fix an  $m \times n$ -bimatrix game  $(A, B)$ . Let  $R$  in  $\mathbb{R}^{m \times n}$  be defined by

$$R_{ij} = \max\{A_{i\ell} \mid \ell \in PB_2(B, e_i)\}.$$

So, rows of  $R$  are constant. Similarly, define  $S$  in  $\mathbb{R}^{m \times n}$  by

$$S_{ij} = \max\{B_{kj} \mid k \in PB_1(A, f_j)\}.$$

The alternative perturbations of  $A$  and  $B$  will be

$$A(\varepsilon_1) = A + \varepsilon_1 R \text{ and } B(\varepsilon_2) = B + \varepsilon_2 S.$$

**Lemma 1** Let  $(A, B)$  be an  $m \times n$ -bimatrix game. A profile  $(p, q)$  is IRE if and only if there exist sequences  $(\varepsilon_1^t)_{t \in \mathbb{N}}$  and  $(\varepsilon_2^t)_{t \in \mathbb{N}}$  of positive real numbers converging to zero, and a sequence  $(p^t, q^t)_{t \in \mathbb{N}}$  in  $\Delta$  converging to  $(p, q)$  such that for all  $t \in \mathbb{N}$

$$(p^t, q^t) \in E(A(\varepsilon_1^t), B(\varepsilon_2^t)).$$

*Proof* Best reply equivalent games have identical equilibrium sets. Since the definition of IRE concerns equilibrium sets of perturbed games, we might as well use other perturbed games as long as they are best reply equivalent. It is easy to verify that  $(A, B)$  and  $(tA, uB)$  are best reply equivalent for any positive real numbers  $t$  and  $u$ , and so are  $(A, B)$  and  $(A + T, B + U)$  if  $T$  is a matrix with constant columns and  $U$  is a matrix with constant rows. Let  $(\varepsilon_1^t)_{t \in \mathbb{N}}$  and  $(\varepsilon_2^t)_{t \in \mathbb{N}}$  be sequences of positive real numbers converging to zero and let  $t \in \mathbb{N}$ . Define  $T$  and  $U$  in  $\mathbb{R}^{M \times N}$  by

$$T_{ij} = \max_{k \in M} A_{kj} \text{ and } U_{ij} = \max_{\ell \in N} B_{i\ell}.$$

Then

$$\begin{aligned} (A(\varepsilon_1^t, \varepsilon_2^t), B(\varepsilon_1^t, \varepsilon_2^t)) &= ((1 - \varepsilon_1^t - \varepsilon_2^t)A \\ &\quad + \varepsilon_1^t R + \varepsilon_2^t T, (1 - \varepsilon_1^t - \varepsilon_2^t)B + \varepsilon_2^t S + \varepsilon_1^t U) \\ &\equiv_b \left( A + \frac{\varepsilon_1^t}{1 - \varepsilon_1^t - \varepsilon_2^t} R, B + \frac{\varepsilon_2^t}{1 - \varepsilon_1^t - \varepsilon_2^t} S \right) \\ &= \left( A \left( \frac{\varepsilon_1^t}{1 - \varepsilon_1^t - \varepsilon_2^t} \right), B \left( \frac{\varepsilon_2^t}{1 - \varepsilon_1^t - \varepsilon_2^t} \right) \right). \end{aligned}$$

Define for all  $t \in \mathbb{N}$  and  $i \in \{1, 2\}$ :  $\varepsilon_i^t = \varepsilon_i^t / (1 - \varepsilon_1^t - \varepsilon_2^t)$ . Then one might as well use the sequences  $(\varepsilon_1^t)_{t \in \mathbb{N}}$  and  $(\varepsilon_2^t)_{t \in \mathbb{N}}$  in combination with perturbed games of the form  $(A + \varepsilon_1^t R, B + \varepsilon_2^t S)$ .  $\square$

*Example 1* Consider the bimatrix game

$$\begin{bmatrix} (1, 1) & (0, 0) \\ (1, 0) & (0, 2) \end{bmatrix}$$

discussed in the Section 1. The game  $(A(\varepsilon_1, \varepsilon_2), B(\varepsilon_1, \varepsilon_2))$  is given by

$$(1 - \varepsilon_1 - \varepsilon_2) \begin{bmatrix} (1, 1) & (0, 0) \\ (1, 0) & (0, 2) \end{bmatrix} + \varepsilon_1 \begin{bmatrix} (1, 1) & (1, 1) \\ (0, 2) & (0, 2) \end{bmatrix} + \varepsilon_2 \begin{bmatrix} (1, 1) & (0, 2) \\ (1, 1) & (0, 2) \end{bmatrix} =$$

$$\begin{bmatrix} (1, 1) & (\varepsilon_1, \varepsilon_1 + 2\varepsilon_2) \\ (1 - \varepsilon_1, 2\varepsilon_1 + \varepsilon_2) & (0, 2) \end{bmatrix},$$

and the alternative perturbation  $(A + \varepsilon_1 R, B + \varepsilon_2 S)$  is

$$\begin{bmatrix} (1 + \varepsilon_1, 1 + \varepsilon_2) & (\varepsilon_1, 2\varepsilon_2) \\ (1, \varepsilon_2) & (0, 2 + 2\varepsilon_2) \end{bmatrix}.$$

In both cases, the top row of the perturbed  $A$ -matrix strictly dominates its bottom row. Hence, in an IRE, player 1 will play top. Player 2 can only respond optimally by choosing left, so the unique IRE is (top, left). Notice that, where the Nash-equilibrium concept does not, the IRE concept treats the following best reply equivalent game differently:

$$\begin{bmatrix} (-1, 1) & (0, 0) \\ (-1, 0) & (0, 2) \end{bmatrix}.$$

Here, the profile (bottom, right) is the unique IRE. Although it heavily relies on the precise situations that are modelled by the two games, we have the opinion that in this example IRE outperforms any Nash-refinement that is invariant under best reply equivalent manipulation.

**Theorem 1** *Let  $(A, B)$  be a bimatrix game. Then  $\text{IRE}(A, B)$  is a non-empty and closed subset of  $E(A, B)$ .*

*Proof* Firstly, we show the non-emptiness. Let  $(\varepsilon_1^t)_{t \in \mathbb{N}}$  and  $(\varepsilon_2^t)_{t \in \mathbb{N}}$  be sequences of positive numbers converging to 0. For all  $t \in \mathbb{N}$ , let  $(p^t, q^t) \in E(A(\varepsilon_1^t), B(\varepsilon_2^t))$ . Due to compactness of the strategy spaces, there exists a subsequence of  $(p^t, q^t)_{t \in \mathbb{N}}$  converging to, say,  $(p, q) \in \Delta$ , which is an element of  $\text{IRE}(A, B)$  by Definition 2 and Lemma 1.

To prove that  $(p, q) \in \text{IRE}(A, B)$  is a Nash equilibrium, we show that  $C(p) \subseteq PB_1(A, q)$  and  $C(q) \subseteq PB_2(B, p)$ . Obviously, it suffices to prove the first statement. Take sequences  $(\varepsilon_1^t)_{t \in \mathbb{N}}$  and  $(\varepsilon_2^t)_{t \in \mathbb{N}}$  of positive numbers converging to 0 and profiles  $(p^t, q^t)$  in  $E(A(\varepsilon_1^t), B(\varepsilon_2^t))$  converging to  $(p, q)$ . Let  $i \in C(p)$ . Then  $i \in C(p^t)$  for  $t$  sufficiently large. Hence, for all  $k \in M$ :

$$e_i A(\varepsilon_1^t) q^t \geq e_k A(\varepsilon_1^t) q^t.$$

Taking  $t$  to infinity, we find for all  $k \in M$ :

$$e_i A q \geq e_k A q.$$

Finally, we show that  $\text{IRE}(A, B)$  is closed. Take a converging sequence  $(p^t, q^t)_{t \in \mathbb{N}}$  in  $\text{IRE}(A, B)$  with limit  $(p, q)$ . For every  $t$ , there are sequences  $(\varepsilon_1^{tk}, \varepsilon_2^{tk})_{k \in \mathbb{N}}$  converging to  $(0, 0)$  and  $(p^{tk}, q^{tk})_{k \in \mathbb{N}}$  converging to  $(p^t, q^t)$  with

$$(p^{tk}, q^{tk}) \in E(A(\varepsilon_1^{tk}), B(\varepsilon_2^{tk})).$$

Consider the sequences  $(\varepsilon_1^{tt}, \varepsilon_2^{tt})_{t \in \mathbb{N}}$  and  $(p^{tt}, q^{tt})_{t \in \mathbb{N}}$ . They demonstrate that  $(p, q)$  is in  $\text{IRE}(A, B)$ .  $\square$

### 3 Strict IRE

Like Robson (1994), we allowed the probabilities of information leakage to be different for the respective players. But requiring them to be equal does not affect the set of informationally robust equilibria. One can go even further: if there is *some* sequence of perturbed games making some profile an IRE, then *any* sequence of perturbed games converging to the original game supports this profile being an

IRE. This section proves the above statement. Firstly, the notion of *strict IRE* is defined, analogously to the way Okada (1984) has refined perfectness to strict perfectness.

**Definition 3** An equilibrium  $(p, q)$  of game  $(A, B)$  is called a *strict IRE* if for all decreasing sequences  $(\varepsilon_1^t, \varepsilon_2^t)_{t \in \mathbb{N}}$  converging to  $(0, 0)$  there is a sequence  $(p^t, q^t)_{t \in \mathbb{N}}$  converging to  $(p, q)$  with  $(p^t, q^t) \in E(A(\varepsilon_1^t, \varepsilon_2^t), B(\varepsilon_1^t, \varepsilon_2^t))$  for all  $t \in \mathbb{N}$ .

**Theorem 2** For any bimatrix game  $(A, B)$  the sets of IRE and strict IRE coincide.

*Proof* Obviously, every strict IRE is an IRE. Conversely, let  $(p, q) \in \text{IRE}(A, B)$  with – see Lemma 1 – associated sequences

$$(\delta_1^t, \delta_2^t) \longrightarrow (0, 0) \text{ and } (p^t, q^t) \longrightarrow (p, q)$$

and  $(p^t, q^t) \in E(A(\delta_1^t), B(\delta_2^t))$  for all  $t \in \mathbb{N}$ . Using subsequences if necessary, we can assume that  $C(p) \subseteq C(p^t) = C(p^{t'})$  and  $C(q) \subseteq C(q^t) = C(q^{t'})$  for all  $t, t' \in \mathbb{N}$ .

Take an arbitrary decreasing sequence  $(\varepsilon_1^t, \varepsilon_2^t)_{t \in \mathbb{N}}$  converging to  $(0, 0)$ . To show that this sequence of perturbations supports  $(p, q)$  as an IRE, we find a  $T \in \mathbb{N}$  and a sequence  $(\hat{p}^t, \hat{q}^t)_{t \geq T}$  converging to  $(p, q)$  with  $(\hat{p}^t, \hat{q}^t) \in E(A(\varepsilon_1^t), B(\varepsilon_2^t))$  for all  $t \geq T$ .

Fix  $T \in \mathbb{N}$  with  $\delta_1^1 > \varepsilon_1^T$  and  $\delta_2^1 > \varepsilon_2^T$ . For every  $t \in \mathbb{N}$ ,  $t \geq T$ , choose  $k(t) \in \mathbb{N}$  such that

$$\delta_1^1 > \varepsilon_1^t > \delta_1^{k(t)} \text{ and } \delta_2^1 > \varepsilon_2^t > \delta_2^{k(t)}.$$

Indeed, for  $i = 1, 2$ , the first inequality  $\delta_i^1 > \varepsilon_i^t$  is automatically fulfilled, since  $\delta_i^1 > \varepsilon_i^T$  and the sequence  $(\varepsilon_1^t, \varepsilon_2^t)_{t \in \mathbb{N}}$  is decreasing. Hence, there are unique  $\lambda(t), \mu(t) \in (0, 1)$  with

$$\varepsilon_1^t = \lambda(t)\delta_1^1 + (1 - \lambda(t))\delta_1^{k(t)} \text{ and } \varepsilon_2^t = \mu(t)\delta_2^1 + (1 - \mu(t))\delta_2^{k(t)}.$$

Define the profile  $(\hat{p}^t, \hat{q}^t)$  by

$$\hat{p}^t = \mu(t)p^1 + (1 - \mu(t))p^{k(t)} \text{ and } \hat{q}^t = \lambda(t)q^1 + (1 - \lambda(t))q^{k(t)}.$$

Since  $(\varepsilon_1^t, \varepsilon_2^t) \longrightarrow (0, 0)$  and  $(\delta_1^t, \delta_2^t) \longrightarrow (0, 0)$ , it follows that  $(\hat{p}^t, \hat{q}^t) \longrightarrow (p, q)$ . It remains to show that  $(\hat{p}^t, \hat{q}^t) \in E(A(\varepsilon_1^t), B(\varepsilon_2^t))$  for all  $t \geq T$ . So let  $t \geq T$ . Because of the similarity, we only show that  $C(\hat{p}^t) \subseteq PB_1(A(\varepsilon_1^t), \hat{q}^t)$ . Take  $i \in C(\hat{p}^t)$ . Because  $C(\hat{p}^t) = C(p^1) = C(p^{k(t)})$  and  $(p^1, q^1)$  is an element of  $E(A(\delta_1^1), B(\delta_2^1))$ , we have for all  $k \in M$ :

$$e_i(A + \delta_1^1 R)q^1 \geq e_k(A + \delta_1^1 R)q^1.$$

Because the rows of  $R$  are constant, we can rewrite this to be

$$e_i A q^1 + \delta_1^1 r_i \geq e_k A q^1 + \delta_1^1 r_k, \quad (1)$$



in which  $r \in \mathbb{R}^m$  is any column of  $R$ . Similarly, for all  $k \in M$ :

$$e_i A q^{k(t)} + \delta_1^{k(t)} r_i \geq e_k A q^{k(t)} + \delta_1^{k(t)} r_k. \quad (2)$$

Adding  $\lambda(t)$  times inequality (1) to  $(1 - \lambda(t))$  times inequality (2) results in ( $k \in M$ )

$$e_i A \hat{q}^t + \varepsilon_1^t r_i \geq e_k A \hat{q}^t + \varepsilon_1^t r_k, \quad (3)$$

which boils down to

$$e_i (A + \varepsilon_1^t R) \hat{q}^t \geq e_k (A + \varepsilon_1^t R) \hat{q}^t \quad (4)$$

for all  $k \in M$ . Hence,  $\hat{p}^t$  is a best response to  $\hat{q}^t$  with respect to the game  $(A + \varepsilon_1^t R, B + \varepsilon_2^t S)$ .  $\square$

Because IRE and strict IRE coincide, Lemma 1 implies that one might as well only look at perturbations of the form  $(A(\varepsilon), B(\varepsilon)) = (A + \varepsilon R, B + \varepsilon S)$ .

**Corollary 1**  $(p, q) \in \text{IRE}(A, B)$  if and only if it is the limit of some trajectory  $(p_\varepsilon, q_\varepsilon)_{\varepsilon \downarrow 0}$  with  $(p_\varepsilon, q_\varepsilon) \in E(A + \varepsilon R, B + \varepsilon S)$ .

#### 4 The structure of IRE

In bimatrix games, the set of Nash equilibria is the union of finitely many Nash components (Jansen 1981). This section shows that the set of informationally robust equilibria of a bimatrix game can be divided into a finite set of components as well. Let  $(A, B)$  be a bimatrix game. By definition, a set of strategy profiles  $G$  is called an IRE component if

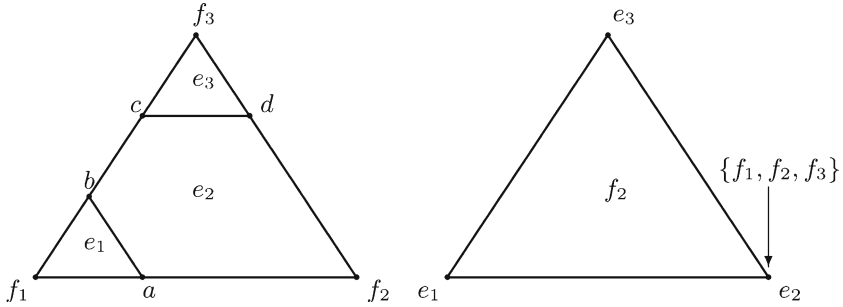
- (i)  $G$  is a convex subset of  $\text{IRE}(A, B)$ ,
- (ii)  $G$  is a product set, i.e.,  $G = G_1 \times G_2$  for some  $G_1 \subseteq \Delta_m, G_2 \subseteq \Delta_n$ ,
- (iii)  $G$  is maximal with respect to properties (i) and (ii).

Replacing  $\text{IRE}(A, B)$  by  $E(A, B)$  yields the definition of a Nash component. To get acquainted with the material, let us start with an example. It shows that different IRE components can be situated in the same Nash component.

*Example 2* Let  $(A, B)$  be

$$\begin{bmatrix} (7, 4) & (2, 5) & (3, 2) \\ (6, 3) & (4, 3) & (5, 3) \\ (4, 2) & (2, 6) & (6, 5) \end{bmatrix}.$$

Figure 2 provides the pure best reply figures. The left-hand side figure, i.e., the one concerning player 1, displays the mixed strategy space of player 2, divided in three parts. Their relative interiors are the areas in which the strategies of player 2 are situated with unique best replies. Four boundary points have been given a name, i.e.,  $a = (2/3)f_1 + (1/3)f_2$ ,  $b = (2/3)f_1 + (1/3)f_3$ ,  $c = (2/3)f_3 + (1/3)f_1$  and  $d = (2/3)f_3 + (1/3)f_2$ . The right-hand side figure shows that  $e_2$  has three pure best replies and all other (mixed) strategies of player 1 have  $f_2$  as their unique best reply. Let  $(p, q)$  be a Nash equilibrium of  $(A, B)$ . If  $PB_2(B, p) = \{f_2\}$ , then  $q$



**Fig. 2** Pure best reply figures of player 1 (left) and player 2 (right)

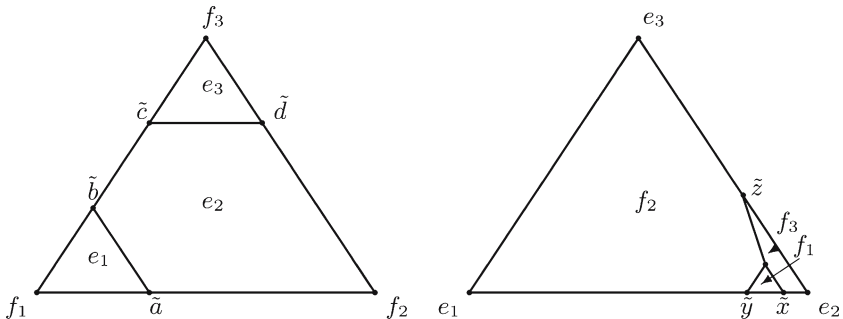
can only be  $f_2$ . The unique best reply on  $f_2$  is  $e_2$ , but  $PB_2(B, e_2) = \{f_1, f_2, f_3\}$ . Hence,  $PB_2(B, p) = \{f_1, f_2, f_3\}$  and  $p$  equals  $e_2$ . Therefore  $e_2$  must be a best reply on  $q$ , so  $q$  is situated in the convex hull of  $a, b, c, d$  and  $f_2$ . We conclude that the unique Nash component is  $\{e_2\} \times \text{conv}(\{a, b, c, d, f_2\})$ .

Let  $\varepsilon$  be a positive number close to zero. The perturbed game  $(A, B, \varepsilon)$  equals

$$\begin{bmatrix} (7 + 2\varepsilon, 4 + 4\varepsilon) & (2 + 2\varepsilon, 5 + 3\varepsilon) & (3 + 2\varepsilon, 2 + 5\varepsilon) \\ (6 + 6\varepsilon, 3 + 4\varepsilon) & (4 + 6\varepsilon, 3 + 3\varepsilon) & (5 + 6\varepsilon, 3 + 5\varepsilon) \\ (4 + 2\varepsilon, 2 + 4\varepsilon) & (2 + 2\varepsilon, 6 + 3\varepsilon) & (6 + 2\varepsilon, 5 + 5\varepsilon) \end{bmatrix}.$$

The pure best reply figures of the perturbed game are depicted in Figure 3. It turns out that the point at which player 2 is indifferent between all his three pure strategies has shifted slightly from  $e_2$  into the interior of the strategy space of player 1. This leads to three areas with a unique pure best reply. Three boundary points have been given a name, i.e.,  $\tilde{x} = (1/2\varepsilon, 1 - 1/2\varepsilon, 0)$ ,  $\tilde{y} = (\varepsilon, 1 - \varepsilon, 0)$  and  $\tilde{z} = (0, 1 - 2\varepsilon, 2\varepsilon)$ . It is easy to infer that the perturbed game has three Nash equilibria:  $(\tilde{x}, \tilde{b})$ , converging to  $(e_2, b)$ ;  $(\tilde{y}, \tilde{a})$ , converging to  $(e_2, a)$ ; and  $(\tilde{z}, \tilde{d})$ , converging to  $(e_2, d)$ .

In the example above, the IRE are all extreme points of the Nash component of the game. The example in the introduction of this paper shows that not all Nash



**Fig. 3** pure best reply figures of the perturbed game  $(A, B, \varepsilon)$

components necessarily contain an IRE. The main result of this section is that for every bimatrix game all IRE components are faces of a Nash component.

**Theorem 3** *Let  $(A, B)$  be a bimatrix game. Then  $\text{IRE}(A, B)$  is the union of finitely many IRE components, each of which is a face of a Nash component, and thereby a polytope.*

In the following proof the phrase “ $(p, q)$  is situated on the face  $F$  of polytope  $P$ ” denotes that  $(p, q)$  is an element of the relative interior of  $F$ . Note that for every  $(p, q) \in P$ , there is exactly one face with this property.

*Proof* The heart of the proof consists of showing the following assertion. Let  $(p, q)$  be an informationally robust equilibrium of the game  $(A, B)$ . Let  $(p', q')$  be situated on the same face of the same component of  $E(A, B)$  as  $(p, q)$ . Then  $(p', q')$  is an element of  $\text{IRE}(A, B)$  as well. Once we have established to show the validity of this assertion, the fact that  $\text{IRE}(A, B)$  is a closed set leads to the observation that IRE components behave like Nash components, which completes the proof.

Hence, let us focus on the assertion above. Because a component is the cartesian product of two polytopes,  $(p', q)$  is situated on the same face as  $(p, q)$  and  $(p', q')$  are. We assume that  $q$  equals  $q'$ , since if we can prove that  $(p', q) \in \text{IRE}$ , we can repeat the argument for  $(p', q')$ , given that  $(p', q) \in \text{IRE}$ . Inside the relative interior of the face of a Nash component the carrier  $C(\cdot)$  and pure best reply correspondence  $PB_2(B, \cdot)$  are constant (see e.g. Jurg 1993, Sect. 2.2). Hence, we have  $C(p) = C(p')$  and  $PB_2(B, p) = PB_2(B, p')$ . Furthermore, since  $(p, q) \in \text{IRE}(A, B)$ , there is a decreasing sequence  $(\varepsilon^t)_{t \in \mathbb{N}}$  with limit 0 and a series of profiles  $(p^t, q^t)_{t \in \mathbb{N}}$  converging to  $(p, q)$  such that  $(p^t, q^t)$  is an equilibrium of the game  $(A(\varepsilon^t), B(\varepsilon^t))$ . For all  $t$ , define

$$\hat{p}^t = p' - p + p^t.$$

Then  $\hat{p}^t$  converges to  $p'$ . For large  $t$ ,  $\hat{p}^t$  is a strategy of player 1, because

$$\sum_{i \in M} \hat{p}_i^t = \sum_{i \in M} p'_i - \sum_{i \in M} p_i + \sum_{i \in M} p_i^t = 1$$

and if  $\hat{p}_i^t < 0$ , then  $i \in C(p) = C(p')$ , so  $p'_i > 0$ . Hence, increasing  $t$  will sufficiently lead to a positive value of  $\hat{p}_i^t$ . The proof is complete when we can show that  $(\hat{p}^t, q^t) \in E(A(\varepsilon^t), B(\varepsilon^t))$ . We have

$$C(\hat{p}^t) \subseteq C(p') \cup C(p) \cup C(p^t) = C(p^t) \subseteq PB_1(A(\varepsilon^t), q^t),$$

so it remains to show that  $C(q^t) \subseteq PB_2(B(\varepsilon^t), \hat{p}^t)$ . Let  $j \in C(q^t)$  and  $\ell \in PB_2(B(\varepsilon^t), \hat{p}^t)$ . Since  $C(q^t) \subseteq PB_2(B(\varepsilon^t), p^t)$ , we have

$$p^t(B + \varepsilon^t S)f_\ell \leq p^t(B + \varepsilon^t S)f_j. \quad (5)$$

Because pure best reply correspondences are upper semi-continuous (see Section 2), for  $t$  sufficiently large we obtain

$$C(q^t) \subseteq PB_2(B(\varepsilon^t), p^t) \subseteq PB_2(B, p) \text{ and } PB_2(B(\varepsilon^t), \hat{p}^t) \subseteq PB_2(B, p').$$

Combining these statements gives

$$\{j, \ell\} \subseteq PB_2(B, p) = PB_2(B, p').$$

This implies that

$$pBf_\ell = pBf_j \text{ and } -p'Bf_\ell = -p'Bf_j. \quad (6)$$

Because the columns of  $S$  are constant, we have

$$p(\varepsilon^t S)f_\ell = p'(\varepsilon^t S)f_\ell \text{ and } p(\varepsilon^t S)f_j = p'(\varepsilon^t S)f_j. \quad (7)$$

Observations (5), (6) and (7) imply

$$(p - p' + p^t)(B + \varepsilon^t S)f_\ell \leq (p - p' + p^t)(B + \varepsilon^t S)f_j.$$

Hence, like  $\ell$ , the strategy  $j$  is an element of  $PB_2(B(\varepsilon^t), \hat{p}^t)$ . We conclude that  $(\hat{p}^t, q^t)$  is an element of  $E(A(\varepsilon^t), B(\varepsilon^t))$ .  $\square$

## 5 Potential games

Potential games have been introduced by Monderer and Shapley (1996). There are many economic situations that can be modeled by potential games. For an overview we refer to Voorneveld (1999). The main virtue of having a potential function for a finite game is that it implies the existence of an (easily traceable) Nash equilibrium in pure strategies. Perhaps the most natural definition of a potential is the *cardinal* (or exact) potential function. On the other hand, the *ordinal* potential generalizes this concept to a much wider class of games and can still be used to obtain the result of this section. Therefore, we give the definition of the latter type of potential.

**Definition 4** A bimatrix game  $(A, B)$  is an ordinal potential game if there exists a function  $P : \Delta \rightarrow \mathbb{R}$  such that for all  $p, p' \in \Delta_m$  and  $q, q' \in \Delta_n$ :

$$\begin{aligned} pAq > p'Aq & \text{ if and only if } P(p, q) > P(p', q), \text{ and} \\ pBq > pBq' & \text{ if and only if } P(p, q) > P(p, q'). \end{aligned}$$

The function  $P$  is called an (ordinal) potential of the game  $(A, B)$ .

It turns out that IRE and the set of strategy pairs at which the potential is maximal always have at least one profile in common.

**Theorem 4** Let  $(A, B)$  be a bimatrix game with ordinal potential  $P$ . Then there exists a pure informationally robust equilibrium that maximizes the potential.

*Proof* Define the  $m \times n$ -matrix  $\bar{P}$  as the restriction of  $P$  to the pure strategy profiles of  $(A, B)$ :

$$\bar{P}_{ij} = P(e_i, f_j). \quad (i \in M, j \in N)$$

By definition of a potential, for all  $i, k \in M$  and all  $j, \ell \in N$ :

$$\begin{aligned} A_{ij} > A_{kj} &\iff \bar{P}_{ij} > \bar{P}_{kj}, \\ B_{ij} > B_{i\ell} &\iff \bar{P}_{ij} > \bar{P}_{i\ell}. \end{aligned} \quad (8)$$

Let us call a matrix satisfying (8) a *potential matrix of*  $(A, B)$ . Firstly, we show that the perturbation  $(A + \varepsilon R, B + \varepsilon S)$  has potential matrix  $\bar{P} + \varepsilon(R + S)$  if  $\varepsilon > 0$  is chosen sufficiently small. Let  $i, k \in M$  and  $j \in N$ . If  $A_{ij} = A_{kj}$ , then  $\bar{P}_{ij} = \bar{P}_{kj}$  and therefore

$$(A + \varepsilon R)_{ij} > (A + \varepsilon R)_{kj} \iff (\bar{P} + \varepsilon R)_{ij} > (\bar{P} + \varepsilon R)_{kj}. \quad (9)$$

If  $A_{ij} > A_{kj}$ , then  $\bar{P}_{ij} > \bar{P}_{kj}$  and we can choose  $\varepsilon$  sufficiently small to obtain the validity of the statements  $(A + \varepsilon R)_{ij} > (A + \varepsilon R)_{kj}$  and  $(\bar{P} + \varepsilon R)_{ij} > (\bar{P} + \varepsilon R)_{kj}$  in (9). Similarly, (9) holds when  $A_{ij} < A_{kj}$  and  $\varepsilon$  is sufficiently small (switch the roles of  $i$  and  $k$ ). Because  $S$  has constant columns we have  $S_{ij} = S_{kj}$ , making (9) equivalent with

$$(A + \varepsilon R)_{ij} > (A + \varepsilon R)_{kj} \iff (\bar{P} + \varepsilon R + \varepsilon S)_{ij} > (\bar{P} + \varepsilon R + \varepsilon S)_{kj}.$$

Similarly, for all  $i \in M$  and all  $j, \ell \in N$  and sufficiently small  $\varepsilon$ :

$$(B + \varepsilon S)_{ij} > (B + \varepsilon S)_{i\ell} \iff (\bar{P} + \varepsilon R + \varepsilon S)_{ij} > (\bar{P} + \varepsilon R + \varepsilon S)_{i\ell}.$$

Hence, the perturbations have potential matrices as well. It is easy to infer that a pure strategy profile maximizing a potential matrix is a Nash equilibrium. There are finitely many pure profiles, so for any sequence of perturbed games converging to  $(A, B)$ , there exists a subsequence of it and a pure profile  $(e_i, f_j)$  such that  $(e_i, f_j)$  is a “potential matrix maximizer” in all games in the subsequence. Since the potential matrices of the perturbed games converge to  $\bar{P}$ ,  $(e_i, f_j)$  is a pure IRE maximizing the potential  $P$ .  $\square$

*Remark 1* A function  $P : \Delta \rightarrow \mathbb{R}$  is called a *cardinal* or *exact* potential of  $(A, B)$  if for all  $p, p' \in \Delta_m$  and all  $q, q' \in \Delta_n$  we have

$$pAq - p'Aq = P(p, q) - P(p', q) \text{ and } pBq - pBq' = P(p, q) - P(p, q').$$

In the case that  $P$  is a cardinal potential,  $P$  is the multilinear extension of  $\bar{P}$ . Along the lines of the proof of Theorem 4 it can be shown that the multilinear extension of  $(\bar{P} + \varepsilon(R + S))$  is a cardinal potential of the perturbed game  $(A + \varepsilon R, B + \varepsilon S)$ .

In general, not all potential maximizers survive. In the following cardinal potential game, the set of potential maximizers is the union of two line segments. IRE selects a single equilibrium.

*Example 3* Consider the game

$$(A, B) = \begin{bmatrix} (1, 1) & (2, 1) \\ (1, 2) & (0, 0) \end{bmatrix}$$

with cardinal potential (matrix)

$$P = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Its set of Nash equilibria  $E(A, B)$  is given by  $([e_1, e_2] \times \{f_1\}) \cup (\{e_1\} \times [f_1, f_2])$ . All equilibria maximize  $P$ . The perturbed game

$$(A(\varepsilon), B(\varepsilon)) = \begin{bmatrix} (1 + 2\varepsilon, 1 + 2\varepsilon) & (2 + 2\varepsilon, 1 + \varepsilon) \\ (1 + \varepsilon, 2 + 2\varepsilon) & (\varepsilon, \varepsilon) \end{bmatrix}$$

has potential matrix

$$P + \varepsilon(R + S) = \begin{bmatrix} 1 + 4\varepsilon & 1 + 3\varepsilon \\ 1 + 3\varepsilon & -1 + 2\varepsilon \end{bmatrix}.$$

Its only Nash equilibrium is  $(e_1, f_1)$ .

## 6 Matrix games

A *zero-sum* or *matrix game* is a bimatrix game  $(A, B)$  with  $B = -A$  and is denoted simply by  $A$ . Recall that in matrix games, the set of Nash equilibria has a product structure, i.e.,  $E(A) = O(A)_1 \times O(A)_2$ , where  $O(A)_1$  are the optimal (maximin) strategies of player 1 and  $O(A)_2$  are the optimal (minimax) strategies of player 2. Since maximin/minimax strategies in combination yield the set of Nash equilibrium profiles, it makes sense to refer to elements of  $O(A)_1$  or  $O(A)_2$  separately as equilibrium strategies. This section shows that also  $\text{IRE}(A)$  has such a product structure. It is, like the Nash equilibrium set, a polytope and an element of it can be found in polynomial time.

Let  $A$  be a zero-sum game. By recalling Figure 1, it is easy to see that outcomes of perturbed games are convex combinations of outcomes of the original game, so each perturbation is a zero-sum game as well. Hence, it suffices to give the perturbations of the payoff to player 1:

$$A(\varepsilon_1, \varepsilon_2)_{ij} = (1 - \varepsilon_1 - \varepsilon_2)A_{ij} + \varepsilon_1 \min_{\ell \in N} A_{i\ell} + \varepsilon_2 \max_{k \in M} A_{kj}.$$

The matrix  $R$  becomes  $((i, j) \in M \times N)$ ,

$$R_{ij} = \max\{A_{i\ell} \mid \ell \in PB_2(-A, e_i)\} = \min_{\ell \in N} A_{i\ell}.$$

Similarly, for  $i$  in  $M$  and  $j$  in  $N$ ,

$$S_{ij} = \max\{-A_{kj} \mid e_k \in PB_1(A, f_j)\} = \min_{k \in M} -A_{kj} = -\max_{k \in M} A_{kj}.$$

By Lemma 1, one might as well consider the perturbed game

$$(A + \varepsilon_1 R, -A + \varepsilon_2 S).$$

This game is best reply equivalent with the zero-sum game

$$A + \varepsilon_1 R - \varepsilon_2 S.$$

Finally, because IRE and strict IRE coincide, one might as well consider the perturbed game

$$A + \varepsilon(R - S).$$

Let  $r \in \mathbb{R}^M$  be any column of  $R$  (they are identical) and let  $s \in \mathbb{R}^N$  be any row of  $S$ .

**Theorem 5** *Let  $A$  be a zero-sum game. Let  $O(A)_1$  and  $O(A)_2$  be the polytopes of optimal strategies of players 1 and 2, respectively. Then  $\text{IRE}(A)$  is a product set, i.e., it can be decomposed:  $\text{IRE}(A) = IO(A)_1 \times IO(A)_2$ .  $IO(A)_1$  is the face of  $O(A)_1$  at which the linear function*

$$O(A)_1 \longrightarrow \mathbb{R}, \quad p \mapsto \langle p, r \rangle$$

*is maximized. Similarly,  $IO(A)_2$  is the face of  $O(A)_2$  at which the linear function*

$$O(A)_2 \longrightarrow \mathbb{R}, \quad q \mapsto \langle -s, q \rangle$$

*is minimized.*

The proof is based on the following idea. We have seen that  $\text{IRE}(A) \subseteq E(A)$ . It appears that the primal concern of a player is to play an optimal strategy of the original game  $A$ . The term  $\varepsilon R$  is of secondary concern to player 1. Hence, he should, within his Nash polytope, maximize this term. The term  $-\varepsilon S$  has no strategic influence to player 1 since the columns of  $S$  are constant. Because of its technical nature the proof has been postponed to the appendix. It requires acquaintance with the Simplex method (e.g. Nemhauser and Wolsey 1988).

The nature of zero-sum games supports the refinement of informationally robustness. For instance, it reduces the harm “not having a poker face” can have, or the disutility that occurs if it is possible to be “cheaten” with small probabilities. Let us give as an illustration a situation in which IRE selects in our opinion the profile that fits best with the context.

**Example 4** Consider a situation in which a penalty shot has been assigned to a soccer team. Let us give the forward taking the penalty three options: aim at the left corner, the right corner, or just give a firm kick. If the forward is skilled, it is obvious that the best thing to do is aim at a corner. If his aiming is poor, however, and he faces an excellent keeper, he would better shoot firmly and hope for the best. The keeper has three pure strategies as well: dive to the left (from the perspective of the forward), dive to the right, or stand still and react on the shot. In our example, depicted in Figure 4, the forward is moderate and we have designed the figures such that he has various optimal strategies. Because the forward cannot aim perfectly, the figures in the matrix do not represent certain outcomes, but expectations.

The keeper has one optimal strategy:  $q = (1/2)(f_1 + f_2)$ . The forward has two extreme optimal strategies:  $p^1 = (1/2)(e_1 + e_2)$  and  $p^2 = (1/6)(e_1 + e_2 + 4e_3)$ . Which one is better? In spite of the fact that  $p^2$  is weakly dominated by  $p^1$ , the concept of IRE recommends  $p^2$ . In the spirit of the concept,  $p^2$  should be played according to the following lines of thought of the forward:

	Keeper		
	dive to the left	dive to the right	stand still
Forward	↓	↓	↓
	aim at left corner:		
	aim at right corner:		
	just kick firmly:		
	$-1$	$1$	$1$
	$1$	$-1$	$1$
	$0$	$0$	$-\frac{1}{2}$

**Fig. 4** A penalty shot

Suppose the keeper can see which corner I am aiming at. Then my chances reduce. On the other hand, if the keeper can see I go for the firm kick, this information is of less value to him.

By using Theorem 5, it is easy to infer that  $(p^2, q)$  is, indeed, the only informationally robust equilibrium; any row  $r$  of  $R$  equals  $(-1, -1, -1/2)$  and  $-1 = \langle r, p^1 \rangle < \langle r, p^2 \rangle = -2/3$ .

## 7 Concluding remarks and future work

Informationally robust equilibria refine Nash equilibria by introducing small probabilities of information leakage. This final section contains brief discussions of the possibilities to generalize information robustness to settings with infinite strategy spaces and more than two players.

Extending the definition of perturbed games to settings with infinite strategy spaces can be done as follows. Assume that player  $i = 1, 2$  has a set  $S_i$  of pure strategies which is a non-empty and compact subset of some finite-dimensional Euclidean space, and that his utility function  $u_i : S_1 \times S_2 \rightarrow \mathbb{R}$  is continuous. With probability  $\varepsilon_1$ , first player's action  $s_1 \in S_1$  is revealed to player two, who will choose a best reply from  $PB_2(s_1) = \operatorname{argmax}_{s_2 \in S_2} u_2(s_1, s_2)$ , a non-empty and compact set. Breaking ties by selecting a best reply that maximizes first player's utility means choosing an element solving  $\max_{s_2 \in PB_2(s_1)} u_1(s_1, s_2)$ , which has a well-defined solution, since  $u_1$  is continuous and  $PB_2(s_1)$  non-empty and compact. Hence, the definition of a perturbed game easily translates to cases with infinite strategy spaces. However, there are relatively simple examples showing that the perturbed utility functions are not continuous. The choice of the tie breaking rule, however, guarantees equilibrium existence in these perturbations (see e.g. Hellwig et al. 1990).

Extending information robustness to games with more than two players requires careful modeling of the timing and content of information leakage. We give three suggestions:

- (i) Each player, but at most one at a time, hears with a small possibility the strategies of all of his opponents. The player can best reply to this observation.
- (ii) Each player, but at most one at a time, reveals with a small possibility his strategy to all of his opponents. The others play an  $(n - 1)$ -person game thereafter.



- (iii) For each ordered pair of players  $(i, j)$ , there is a slight chance that  $i$  finds out the action of player  $j$ .

The proper model of information leakage in such larger games probably depends on the exact situation being studied and is an interesting direction for further research.

## Appendix

In order to prove Theorem 5, a result is needed from Linear Algebra, providing sufficient conditions for convergence of solution sets of perturbed systems of linear equations.

*Claim* Let  $D$  be an  $m \times n$ -matrix and let  $(d^t)_{t \in \mathbb{N}}$  be a sequence in  $\mathbb{R}^m$  converging to  $d$ . Let for all  $t$  in  $\mathbb{N}$ ,  $F^t \subset \mathbb{R}^n$  be the set of feasible points of the system of equations  $\{x \in \mathbb{R}_+^n \mid Dx = d^t\}$ . Let  $F$  be the set of feasible points of  $\{x \in \mathbb{R}_+^n \mid Dx = d\}$ . Suppose there exists a uniform bound  $M \in \mathbb{N}$ , i.e.,  $\|x\| \leq M$  for all  $x \in \bigcup F^t$ . If all solution sets  $F^t$  are non-empty, then  $F^t$  converges to  $F$  in the sense that

- (i) if  $\hat{x}^t \in F^t$  for all  $t \in \mathbb{N}$  and  $\lim_{t \rightarrow \infty} \hat{x}^t = \hat{x}$ , then  $\hat{x} \in F$ ,
- (ii) for all  $\hat{x} \in F$  there exists a sequence  $(\hat{x}^t)_{t \in \mathbb{N}}$  in  $(F^t)_{t \in \mathbb{N}}$  converging to  $\hat{x}$ .

*Proof* It is easy to infer statement (i) by a continuity argument. The difficult part is to show that any element of  $F$  is the limit of some sequence in  $(F^t)_{t \in \mathbb{N}}$ . The proof will be by induction to  $n$ ; the number of columns. The case  $n = 1$  is left to the reader. We distinguish between two cases.

*Case I:* There exists a strictly positive element  $s \in \mathbb{R}_{++}^n$  of  $F$ .

Linear operations like adding rows to others, or multiplying a row with a non-zero number will not change the solutions sets, nor the feature that the constraint vectors converge. Hence, without loss of generality,  $D$  has an echelon form

$$D = \begin{bmatrix} I_r & M \\ \bar{0} & \bar{0} \end{bmatrix},$$

in which  $r$  is the rank of  $D$ ,  $I_r$  is an identity matrix,  $M$  is some matrix with  $r$  rows and  $n - r$  columns and the zeros represent zero matrices. Because  $F^t \neq \emptyset$  for all  $t \in \mathbb{N}$ , we have that  $d_i = d_i^t = 0$  for all  $t \in \mathbb{N}$  and all  $i > r$ . Hence, we might as well remove the  $m - r$  zero-rows of  $D$ , which boils down to assuming that  $D$  is of full rank:  $r = m$ . Let  $q = (d_1, \dots, d_m, 0, \dots, 0) \in \mathbb{R}^n$ . Then  $Dq = d$ . Similarly, let  $q^t = (d^t, \bar{0}) \in \mathbb{R}^n$ , so  $Dq^t = d^t$ . Define  $s^t = s + q^t - q$ . Let  $\delta > 0$  be such that  $s_i > \delta$  for all  $i \leq n$ . Then  $s_i^t > (1/2)\delta$  for large  $t$  and  $i \leq n$ .

Let  $\hat{x}$  be any element of  $F$ . Define  $\hat{x}^t = \hat{x} + q^t - q$ . Then  $D\hat{x}^t = d^t$  and  $\hat{x}^t \rightarrow \hat{x}$ . Let  $\lambda^t = \min\{\lambda \in [0, 1] \mid \lambda s^t + (1 - \lambda)\hat{x}^t \geq \bar{0}\}$  and define  $\tilde{x}^t = \lambda^t s^t + (1 - \lambda^t)\hat{x}^t \in F^t$ . Let  $\varepsilon > 0$ . Choose  $t$  so large that  $\hat{x}_i^t > -\varepsilon$  for all  $i$ . If  $\hat{x}^t \notin \mathbb{R}_+^n$ , then

$$\lambda^t = \max_{i \leq n} \frac{-\hat{x}_i^t}{s_i^t - \hat{x}_i^t} \leq \frac{\varepsilon}{\frac{1}{2}\delta}.$$

Since  $\delta$  is fixed and  $\varepsilon$  can be chosen to be as small as desired,  $\lambda^t$  tends to 0. Hence,  $\|\hat{x} - \tilde{x}^t\|$  converges to 0 if  $t$  tends to  $\infty$ . This ends Case I.

*Case II:* For some  $i \leq n$ ,  $x_i = 0$  for all  $x \in F$ .

Without loss of generality, choose  $i = n$ . Let  $\delta^t = \min_{x \in F^t} x_n$ . Let  $\delta$  be an accumulation point of  $(\delta^t)_{t \in \mathbb{N}}$ . Because of the uniform bound  $M$ , there exists a sequence  $x^{t_1}, x^{t_2}, x^{t_3}, \dots$  converging to, say,  $x$  with  $x_n^{t_k} = \delta^{t_k}$  and  $\lim_{k \rightarrow \infty} \delta^{t_k} = \delta$ . By continuity,  $x \in F$  and  $\lim_{k \rightarrow \infty} \delta^{t_k} = x_n = 0$ . Hence, 0 is the only accumulation point;  $\lim_{t \rightarrow \infty} \delta^t = 0$ . Substitute, for all  $t \in \mathbb{N}$ ,  $x_n = \delta^t$  in the equation set  $Dx = d^t$ . The solution sets may become smaller, but remain non-empty. By now, the right column can be removed from all sets of equations and we obtain a setting with one-dimension less. Hence, we can apply the induction hypothesis. For an arbitrary element  $\hat{x} = (\hat{x}_1, \dots, \hat{x}_{n-1}, 0)$  of  $F$ , we can give an element  $(\hat{x}_1^t, \dots, \hat{x}_{n-1}^t, \delta^t)$  in  $F^t$  close to  $\hat{x}$ .  $\square$

Notice that if the constraint matrix  $D$  is perturbed as well, convergence is not guaranteed. For example, if

$$D^t = \begin{bmatrix} 1 - \frac{1}{t} & 1 + \frac{1}{t} \\ 1 + \frac{1}{t} & 1 - \frac{1}{t} \end{bmatrix} \text{ and } d^t = \begin{bmatrix} 2 \\ 2 \end{bmatrix},$$

the solution sets of  $F^t$  all equal  $\{(1, 1)\}$ , while the solution set of  $F$  equals  $\{(x, 2 - x) : x \in [0, 2]\}$ .

*Proof of Theorem 5* Because for every  $\varepsilon > 0$ ,  $E(A + \varepsilon(R - S))$  is a product set and a polytope and because  $IRE(A)$  coincides with strict  $IRE(A)$  (Theorem 2),  $IRE(A)$  is a product set and a polytope as well, say  $IRE(A) = IO(A)_1 \times IO(A)_2 \subseteq \Delta_m \times \Delta_n$ . The assertions concerning  $IO(A)_1$  and  $IO(A)_2$  are so similar that we suffice with the proof of the latter. Assume without loss of generality that  $A > 0$ . Then  $R$  is as well a strictly positive matrix and  $S$  is a strictly negative matrix. Furthermore,  $v(A)$ , the value of the game, is strictly positive. Let  $(\varepsilon^t)_{t \in \mathbb{N}}$  be a decreasing sequence with limit 0.  $O(A + \varepsilon^t(R - S))_2$  is the set of optimal solutions of the linear program

$$\begin{array}{ll} \text{minimize } v & \text{subject to} \\ v \in \mathbb{R}_+, q \in \mathbb{R}_+^N & \end{array} \quad \begin{bmatrix} 0 & 1 & \dots & \dots & 1 \\ 1 \\ \vdots & -A + \varepsilon^t(S - R) \\ 1 \end{bmatrix} \begin{bmatrix} v \\ q_1 \\ \vdots \\ q_n \end{bmatrix} \geq \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

The left column will be referred to as column  $v$ , the top row as row 0 and each other row by its corresponding pure strategy: row  $i$  ( $i \in M$ ). If we would like to apply the Simplex method, for each row a slack variable has to be added, except for row 0, since  $\sum_{j \in N} q_j$  has to equal 1. We get

$$\begin{array}{ll} \text{minimize } \left\langle e_v, \begin{bmatrix} v \\ q \\ p \end{bmatrix} \right\rangle & \text{s.t.} \\ v \in \mathbb{R}_+, q \in \mathbb{R}_+^N, p \in \mathbb{R}_+^M & \end{array} \quad \begin{bmatrix} 0 & 1 & \dots & \dots & 1 & 0 & \dots & 0 \\ 1 \\ \vdots & -A + \varepsilon^t(S - R) & & -I_m \\ 1 \end{bmatrix} \begin{bmatrix} v \\ q \\ p \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Here,  $e_v$  denotes the unit vector of  $\mathbb{R} \times \mathbb{R}^N \times \mathbb{R}^M$  corresponding to  $v$  and  $I_m$  denotes the  $m \times m$  identity matrix. By adding row 0 of the table  $\varepsilon^t r_i$  times to row  $i$  ( $i \in M$ ), the table becomes independent of the matrix  $R$ , except for the constraint vector. The resulting table will be denoted by  $LP^t$ :

$$\text{minimize } \left\langle e_v, \begin{bmatrix} v \\ q \\ p \end{bmatrix} \right\rangle \text{ s.t. } \begin{bmatrix} 0 & 1 & \cdots & 1 & 0 & \cdots & 0 \\ 1 & & & & & & \\ \vdots & -A + \varepsilon^t S & & -I_m & & & \\ 1 & & & & & & \end{bmatrix} \begin{bmatrix} v \\ q \\ p \end{bmatrix} = \begin{bmatrix} 1 \\ \varepsilon^t r_1 \\ \vdots \\ \varepsilon^t r_m \end{bmatrix}. \quad (10)$$

Denote the constraint matrix in the program  $LP^t$  by  $D^t$ . The program and constraint matrix obtained by substituting  $\varepsilon^t = 0$  will be called  $LP$  and  $D$ , respectively. They correspond to the non-perturbed game  $A$ . After having performed the Simplex method, the table has become of the following form:<sup>2</sup>

$$\text{minimize } \langle a^t, [v, q, p]^T \rangle \text{ s.t. } B^t [v, q, p]^T = b^t. \quad (11)$$

$v, q, p \geq 0$

Let us recall the features of the Simplex method that are important for our purpose. The final object vector  $a^t \in \mathbb{R}_+ \times \mathbb{R}_+^N \times \mathbb{R}_+^M$  is non-negative and equals the sum of the original object vector  $e_v$  and some linear combination of the rows of  $LP^t$ . The main principle of the Simplex method is, that one might as well optimize the final object vector, because for any row  $D_{i,\cdot}^t$ , the inner product  $\langle D_{i,\cdot}^t, x \rangle$  is independent on  $x$  (as long as  $x$  is chosen feasible). The set of optimal points consists of all feasible points with inner product zero with the final object vector. Because the tables consists of linear equations, we can normalize them such that for all  $t \in \mathbb{N}$ , all numbers in  $B^t$ ,  $b^t$  and  $a^t$  are in some compact segment, e.g.,  $[-1, 1]$ . Hence, by taking a suitable subsequence of the sequence  $(\varepsilon^t)_{t \in \mathbb{N}}$ , we can accomplish that  $B^t$ ,  $b^t$  and  $a^t$  converge to, say,  $B$ ,  $b$  and  $a$ , respectively. This limit (minimize  $\langle a, x \rangle$  s.t.  $Bx = b$ ) is a table for the original game and could have been obtained by applying the Simplex method on  $LP$ . Hence,  $a$  equals  $e_v$  plus some linear combination of the rows of  $LP$ :

$$a = e_v + \sum_{i=0}^m c_i D_i. \quad \text{for some } c \in \mathbb{R} \times \mathbb{R}^M. \quad (12)$$

Because  $v(A)$  is strictly positive, we have that  $x_v = v(A) > 0$  for all optimal points, so  $a_v = 0$ . Focussing at the first column of  $LP$ , equation (12) gives

$$0 = a_v = (e_v)_v + \sum_{i=0}^m c_i D_{iv} = 1 + \sum_{i=1}^m c_i. \quad (13)$$

We have that  $a_i^t > 0$  for large  $t$  and all  $i \in C(a)$ . Hence, all variables corresponding to elements of  $C(a)$  have value 0 in any optimal point and all corresponding columns can be removed<sup>3</sup> from the tables  $LP^t$  and  $LP$  without changing optimal sets: columns in  $C(a) \cap M$  correspond to pure strategies on which player 1 can

<sup>2</sup> The  $\top$  denotes that the vector is transposed.

<sup>3</sup> The removed variables of course still have to be stored and are set to be zero.

put some weight while playing optimal in the original game  $A$  and columns in  $C(a) \cap N$  correspond to pure strategies on which player 2 does not put positive weight in any equilibrium of  $A$ . Denote the complement of the carrier of  $a$  by  $Z(a)$  (the ‘zero part’ of  $a$ ):

$$Z(a) = \{i : a_i = 0\}.$$

Denote the matrices  $D^t$  and  $D$  of which the redundant columns have been deleted by  $\bar{D}^t$  and  $\bar{D}$ , respectively. Similarly, let  $\bar{e}_v = (1, 0, \dots, 0) \in \mathbb{R}^{Z(a)}$  be the first unit vector of  $\mathbb{R}^{Z(a)}$ , let  $\bar{s} \in \mathbb{R}^{Z(a)}$  be the restriction of the vector  $(0, s, 0, \dots, 0) \in \mathbb{R} \times \mathbb{R}^N \times \mathbb{R}^M$  and let  $\bar{a}^t$  be the restriction of  $a^t$  (so  $\bar{a}^t = \bar{0}$  for all  $t$ ). We can omit these columns as well in equation (12):

$$\bar{0} = \bar{a} = \bar{e}_v + \sum_{i=0}^m c_i \bar{D}_i.$$

Adding the rows of  $\bar{D}^t$  to  $\bar{e}_v$ , weighted by the same combination  $c$ , results in

$$\bar{e}_v + \sum_{i=0}^m c_i \bar{D}_i^t = \sum_{i=0}^m c_i (\bar{D}_i^t - \bar{D}_i) = \sum_{i=1}^m c_i \varepsilon^t \bar{s} = (-1) \cdot \varepsilon^t \bar{s}.$$

To infer the second equality, consider program  $LP^t$ , (equation (10)): the difference between row  $i$  of  $LP^t$  and row  $i$  of  $LP$  is  $\varepsilon^t$  times the vector  $(0, s, 0, \dots, 0)$  for all  $i \geq 1$ . For the third equality we refer to (13). Hence, for all  $t \in \mathbb{N}$ , instead of minimizing  $\langle \bar{e}_v, x \rangle$ , we might as well minimize  $\langle -\varepsilon^t \bar{s}, x \rangle$ , or  $\langle -\bar{s}, x \rangle$ . Call the alternative optimization problem  $ALP^t$ :

$$\begin{array}{ll} \text{minimize } \langle -\bar{s}, x \rangle & \text{s.t. } \bar{D}^t x = [1, \varepsilon^t r_1, \dots, \varepsilon^t r_m]^\top. \\ & x \in \mathbb{R}_+^{Z(a)} \end{array}$$

Let us repeat the results so far. For all  $t \in \mathbb{N}$ , the set  $O(A + \varepsilon^t(R - S))_2$  is described by  $ALP^t$  in the sense that for all  $q \in \Delta_n$ , the following statements are equivalent:

- (i)  $q \in O(A + \varepsilon^t(R - S))_2$  and
- (ii)  $q_j = 0$  for all  $j \in N \cap C(a)$  and  $q_j = x_j$  for all  $j \in N \cap Z(a)$  and some optimal solution  $x \in \mathbb{R}_+^{Z(a)}$  of  $ALP^t$ .

Consequently, the program obtained by substituting  $\varepsilon^t = 0$  in  $ALP^t$  will be called  $ALP$ . The set of feasible points of  $ALP$  corresponds to  $O(A)_2$  in the sense that for all  $q \in \Delta_n$ :  $q \in O(A)_2$  if and only if  $q_j = 0$  for all  $j \in N \cap C(a)$  and  $q_j = x_j$  for all  $j \in N \cap Z(a)$  and some *feasible* point  $x \in \mathbb{R}_+^{Z(a)}$  of  $ALP$ . The optimal set of  $ALP$  corresponds to the face of  $O(A)_2$  of which Theorem 5 claims that it coincides with  $IO(A)_2$ . Hence, we are done if we can show that the optimal set of  $ALP^t$  converges to the optimal set of  $ALP$ .

After having performed the Simplex method on table  $ALP^t$ , we get again a table of the form

$$\begin{array}{ll} \text{minimize } \langle h^t, x \rangle & \text{s.t. } G^t x = g^t. \\ & x \geq 0 \end{array}$$

Here,  $h^t \in \mathbb{R}_+^{Z(a)}$ . The following lines of argumentation copy the one just after equation (11), so details are omitted. Assume that  $h^t$  converges to  $h$ . Then

$$h = -\bar{s} + \sum_{i=0}^m \bar{c}_i \bar{D}_i. \quad \text{for some } \bar{c} \in \mathbb{R} \times \mathbb{R}^M. \quad (14)$$

We have that  $h_i^t > 0$  for large  $t$  and all  $i \in C(h)$ . Columns corresponding to elements of  $C(h)$  are removed from the tables  $ALP^t$  and  $ALP$  without changing optimal sets. Denote the matrices  $\bar{D}^t$  and  $\bar{D}$  of which the redundant columns have been deleted by  $\hat{D}^t$  and  $\hat{D}$ , respectively. Similarly, let  $\hat{e}_v \in \mathbb{R}^{Z(h)}$  be the first unit vector of  $\mathbb{R}^{Z(h)}$ , let  $\hat{s} \in \mathbb{R}^{Z(h)}$  be the restriction of  $\bar{s}$ . Omit the redundant columns in equation (14):

$$\bar{0} = -\hat{s} + \sum_{i=0}^m \bar{c}_i \hat{D}_i.$$

If we add the rows of  $\hat{D}^t$  to  $-\hat{s}$ , weighted by combination  $\bar{c}$ , we obtain

$$-\hat{s} + \sum_{i=0}^m \bar{c}_i \hat{D}_i^t = \sum_{i=0}^m \bar{c}_i (\hat{D}_i^t - \hat{D}_i) = \sum_{i=1}^m \bar{c}_i \varepsilon^t \hat{s}.$$

The object vector  $-\hat{s}$  manifests to be a linear combination of the rows of  $\hat{D}^t$ . Hence, the linear function  $\langle -\hat{s}, \cdot \rangle$  is constant on the polytope  $F^t = \{x \in \mathbb{R}_+^{Z(h)} \mid \hat{D}^t x = [1, \varepsilon^t r]^\top\}$ , say  $k^t = \langle -\hat{s}, x \rangle$  for all  $x \in F^t$ . Add to all rows of  $ALP^t$  but the first, the equation  $\varepsilon^t \langle -s, x \rangle = \varepsilon^t k^t$ , to obtain

$$F^t = \left\{ x \in \mathbb{R}_+^{Z(h)} \mid \hat{D}x = [1, \varepsilon^t(r_1 + k^t), \dots, \varepsilon^t(r_m + k^t)]^\top \right\}.$$

Observe that the constraint matrix of this description is no longer dependent on  $t$ . Conclusion: we have found a description of the form  $(\hat{D}x = d^t, x \geq 0)$  of the optimal set of  $ALP^t$  and a description  $(\hat{D}x = (1, 0, \dots, 0), x \geq 0)$  of the optimal set of  $ALP$ . Apply Claim 7 and conclude the validity of Theorem 5.  $\square$

## References

- Bagwell, K.: Commitment and observability in games. *Games Econ Behav* **8**, 271–280 (1995)
- Hellwig, M., Leininger, W., Reny, P., Robson, A.: Subgame perfect equilibrium in continuous games of perfect information: an elementary approach to existence and approximation by discrete games. *J Econ Theor* **52**, 406–422 (1990)
- Jansen, M.: Maximal nash subsets for bimatrix games. *Naval Res Logist Q* **28**, 85–101 (1981)
- Jurg, A.: Some topics in the theory of bimatrix games. Ph. D. thesis, University of Nijmegen (1993)
- Kohlberg, E., Mertens, J.F.: On strategic stability of equilibria. *Econometrica* **54**, 1003–1037 (1986)
- Matsui, A.: Information leakage forces cooperation. *Games Econ Behav* **1**, 94–115 (1989)
- Monderer, D., Shapley, L.: Potential games. *Games Econ Behav* **14**, 124–143 (1996)
- Myerson, R.B.: Refinements of the nash equilibrium point concept. *Int J Game Theor* **7**, 73–80 (1978)

- Nemhauser, G., Wolsey, L.: Integer and combinatorial optimization. New York: Wiley 1988
- Okada, A.: Strictly perfect equilibrium points in strategic games. *Int J Game Theor* **13**, 145–154 (1984)
- Reny, P., Robson, A.: Reinterpreting mixed strategy equilibria: a unification of the classical and bayesian views. *Games Econ Behav* **48**, 355–384 (2004)
- Robson, A.: An ‘informationally robust equilibrium’ for two-person nonzero sum games. *Games Econ Behav* **7**, 233–245 (1994)
- Selten, R.: Reexamination of the perfectness concept for equilibrium points in extensive games. *Int J Game Theor* **4**, 25–55 (1975)
- Solan, E., Yariv, L.: Games with espionage. *Games Econ Behav* **47**, 172–199 (2004)
- Van Damme, E.: Stability and perfection of Nash equilibria. Berlin Heidelberg New York: Springer 1991
- Voorneveld, M.: Potential games and interactive decisions with multiple criteria. Ph. D. thesis, Tilburg University, Center dissertation series (1999)